

(19)



JAPANESE PATENT OFFICE

PATENT ABSTRACTS OF JAPAN

(11) Publication number: **05249987 A**(43) Date of publication of application: **28.09.93**

(51) Int. Cl.

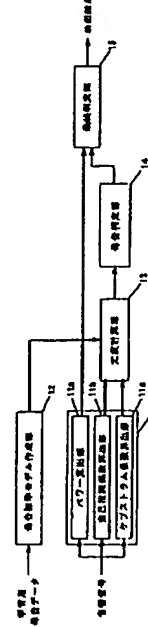
G10L 3/00
G10L 9/00
(21) Application number: **04050327**(22) Date of filing: **09.03.92**(71) Applicant: **MATSUSHITA ELECTRIC IND CO LTD**(72) Inventor: **NAKATO YOSHIHISA**
NORIMATSU TAKESHI(54) **VOICE DETECTING METHOD AND DEVICE**

COPYRIGHT: (C)1993,JPO&Japio

(57) Abstract:

PURPOSE: To automatically detect voice with high precision with a comparatively simple structure in a voice detecting device for detecting only the voice, which is used as the pretreatment of a voice recognizing device.

CONSTITUTION: A voice detecting device has a tolerance calculating part 13 for extracting a plurality of characteristic quantities from input signal every fixed time by a characteristic extracting part 11 and calculating the logarithmic tolerance with a vowel standard model formed by use of a number of learning data of vowel by a vowel standard model forming part 12: and a vowel judging part 14 for calculating a frame average logarithmic tolerance by collectively using the logarithmic tolerances for several frames and detecting vowels by comparison with a proper threshold. According to the number judged as vowels by the vowel judging part, it is judged by a final judging part 15 whether this section is a voice or not.



(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)特許出願公開番号

特開平5-249987

(43)公開日 平成5年(1993)9月28日

(51)Int.Cl.⁵

G10L 3/00
9/00

識別記号

513 B 8842-5H
A 8946-5H

庁内整理番号

FI

技術表示箇所

審査請求 未請求 請求項の数3(全7頁)

(21)出願番号 特願平4-50327

(22)出願日 平成4年(1992)3月9日

(71)出願人 000005821

松下電器産業株式会社
大阪府門真市大字門真1006番地

(72)発明者 中藤 良久

大阪府門真市大字門真1006番地 松下電器
産業株式会社内

(72)発明者 則松 武志

大阪府門真市大字門真1006番地 松下電器
産業株式会社内

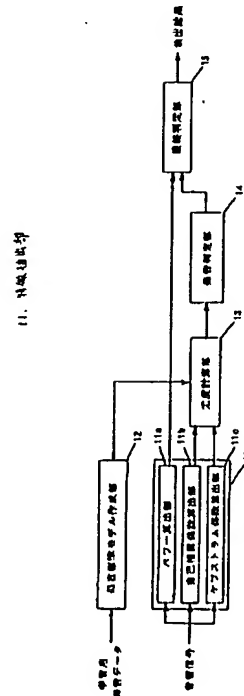
(74)代理人 弁理士 小鍛冶 明 (外2名)

(54)【発明の名称】 音声検出方法および音声検出装置

(57)【要約】

【目的】 本発明は音声認識装置等の前処理として用いられる、音声のみを検出する音声検出装置に関するもので、比較的簡単な構成で自動的にしかも高精度に音声の検出を行うことができる音声検出装置を提供することを目的とする。

【構成】 入力信号から特徴抽出部11にて一定時間毎に複数の特徴量を抽出し、母音標準モデル作成部12であらかじめ多数の母音の学習データを用いて作成した母音標準モデルとの対数尤度を計算する尤度計算部13と、数フレーム分の対数尤度を一括して用いてフレーム平均対数尤度を計算し、適当な閾値と比較することで母音を検出する母音判定部14と、母音判定部により母音と判定された個数により、最終判定部15にてその区間が音声か否かを判定する。



【特許請求の範囲】

【請求項1】 入力信号からフレーム単位（一定時間毎）に抽出した音声の特徴付ける1次以上の自己相関係数もしくは偏自己相関係数と1次以上のケプストラム係数もしくはメルケプストラム係数のうち少なくとも1つの特徴量を用いて、その数フレーム分を一括して用いることにより母音の存在確率を求め母音検出を行い、音声のみを検出することを特徴とする音声検出方法。

【請求項2】 あらかじめ多数の母音の学習データについてフレーム単位に抽出した音声の特徴付ける1次以上の自己相関係数もしくは偏自己相関係数と1次以上のケプストラム係数もしくはメルケプストラム係数のうち少なくとも1つの特徴量を用いて、数フレーム分を一括して用いて母音毎に標準モデルの作成を行い、前記母音標準モデルを用いて音声を検出することを特徴とする音声検出方法。

【請求項3】 入力信号から一定時間毎に音声の特徴付ける1次以上の自己相関係数もしくは偏自己相関係数と1次以上のケプストラム係数もしくはメルケプストラム係数を抽出する特徴抽出部と、あらかじめ多数の母音の学習データについて前記特徴抽出部で抽出した特徴量を用いて母音毎の平均値と共分散行列を算出し、母音毎の標準モデルを作成する母音標準モデル作成部と、入力信号からフレーム単位に前記特徴抽出部で抽出した1次以上の自己相関係数もしくは偏自己相関係数と1次のケプストラム係数もしくはメルケプストラム係数のうち少なくとも1つ以上の特徴量について、前記母音標準モデル作成部にて作成した各母音標準モデルとの対数尤度を計算する尤度計算部と、母音検出しようとするフレームとその前後数フレームにおいて前記尤度計算部にて計算された対数尤度を用いて各母音毎にフレーム平均対数尤度を計算し、ある適当な閾値と比較することで母音かそれ以外かを判定する母音判定部と、パワーの一定レベル以上の入力信号の塊について前記母音判定部により、いずれかの母音と判定された母音サンプルの個数がある適当な閾値以上のときにその塊を音声と判定する最終判定部とを備えたことを特徴とする音声検出装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、非定常雑音の存在する実環境下において、音声認識装置の前処理等で使われる、音声のみを検出する音声検出方法および音声検出装置に関する。

【0002】

【従来の技術】 音声認識等の音声処理を行う装置では、音声以外の非定常雑音が入力され誤って音声と判断されると誤認識を生じる。そこで、入力された信号が正確に音声であるかどうかを判定できる音声検出装置が必要とされる。

【0003】 従来の音声検出装置では、処理の簡素化の

ための入力信号のパワー値が閾値よりも大きい部分を音声と判断する方法が一般的に行われる。しかし音声認識の行われる実環境で使用することを考えると、紙などの資料をめくる音や、息吹きなどのマイクローフンの振動によって起こるノイズ、あるいは動物の鳴き声等の音声以外のパワーの大きな様々な音が入力される可能性があり、パワーだけでは音声の検出はできない。

【0004】 そこで、パワー以外の複数の音声の特徴量を用いて入力信号が音声であるか非音声であるかの判定をする方法が幾つか提案されている。例えば、「実環境下での音声／非音声の判別」（石田明・小畑秀文、日本音響学会誌7巻12号（1991））による方法がある。これは、日常の実験室やオフィスなどで発生する種々の非定常雑音と音声とを区別するのに有効な音響的特徴量を用いて、実環境下での音声／非音声の判別を行っている。具体的には、音声の中のパワーの大きい部分において、母音と見なせる部分がどの程度存在するかによって、音声／非音声の判別を行っており、用いる音響的特徴量としては、

- (a) 周期性
- (b) ピッチ周波数
- (c) 最適線形予測次数
- (d) 5母音との距離
- (e) ホルマンントの鋭さ

の5種類の特徴量を求め、各特徴量毎に上限値あるいは下限値を決定し、その大小関係により音声と非音声を判別する。

【0005】

【発明が解決しようとする課題】 しかしながら上記の音声／非音声判別装置では、音声の中の各母音の特徴に基づいた特徴量は使用されておらず、音声の各母音の検出に適した母音毎の標準モデルを用いる方法による高精度な音声検出方式が必要とされる。また、各特徴量毎に上限値あるいは下限値を決定し、その大小関係により音声／非音声の判別を行う方法では、特徴量の数が増えた場合特徴量毎に閾値を設定することが困難であると同時に、定常雑音が付加された場合などのような特徴量の値の変動に対して頑健であるとは言えない。さらに、音声、特に母音は1分析フレーム毎に母音性を判定するより、数フレーム分を1塊に考えて判定する方がより信頼性がある判定法であるといえる。

【0006】 本発明は、上記の課題を解決するもので、音声認識等の音声信号処理に適した高性能な音声検出装置を提供することを目的とする。本発明は、音声の各母音の検出に適した母音毎の標準モデルを用いることで、音声の各母音の検出に基づいた音声検出装置を提供する。さらに、音声であるかそれ以外であることを表した特徴量を総合的に判定するための評価値として、母音性を判定するのに有用と考えられる数フレーム分を1塊に考えて算出されるフレーム平均対数尤度を用いることで、

閾値の設定が比較的容易であると同時に、定常雑音が付加された場合などのような特徴量の値の変動に対してある程度頑健性を持った、比較的簡単で高性能な音声検出装置を提供することを目的とする。

【0007】

【課題を解決するための手段】本発明は上記課題を解決するために、入力信号からフレーム単位（一定時間毎）に抽出した音声の特徴付ける1次以上の自己相関係数と1次以上のケプストラム係数のうち少なくとも1つの特徴量を用いて、その数フレーム分を一括して用いることにより母音の存在確率を求めて母音検出を行い、これにより音声のみを検出することを特徴とするものである。

【0008】また、本発明は、あらかじめ多数の母音の学習データについてフレーム単位に抽出した音声の特徴付ける1次以上の自己相関係数と1次以上のケプストラム係数のうち少なくとも1つの特徴量を用いて、数フレーム分を一括して用いて母音毎に標準モデルの作成を行い、前記母音標準モデルを用いて音声を検出することを特徴とするものである。

【0009】さらに、本発明の音声検出装置は、母音を検出することを主眼として、入力信号の一定時間毎の1次以上の自己相関係数、1次以上のケプストラム係数等の複数の音声の特徴量を抽出する特徴量抽出部と、あらかじめ多数の母音の学習データについて前記特徴抽出部で抽出した特徴量を用いて母音毎の平均値と共分散行列を算出し、母音毎の標準モデルを作成する母音標準モデル作成部と、入力信号からフレーム単位に前記特徴抽出部で抽出した1次以上の自己相関係数と1次のケプストラム係数のうち少なくとも1つの特徴量について、前記母音標準モデル作成部にて作成した各母音標準モデルとの対数尤度を計算する尤度計算部と、前記尤度計算部にて計算された前後数フレームについて各母音毎にフレーム平均対数尤度を計算し、ある適当な閾値とを比較することで母音がそれ以外かを判定する母音判定部と、パワーの一定レベル以上の入力信号の塊について前記母音判定部によりいずれかの母音と判定された母音サンプルの個数の割合がある適当なしきい値以上のときにその塊を音声と判定する最終判定部とを備えたものである。

【0010】

【作用】本発明は、上記した構成により、音声中の各音韻の特徴に基づく母音検出に適した特徴量を用い、あらかじめ信頼性の高い多数の母音データを用いて母音毎に母音標準モデルを作成し、数フレーム分を一括して統計的手法により母音の検出を行い、音声のみを検出することで、高性能な音声検出が可能となる。

【0011】

【実施例】以下本発明の一実施例について説明する。

（図1）は本発明の一実施例の全体構成を示すブロック構成図である。（図1）において、11は音声検出のための複数の特徴量を抽出する特徴抽出部で、1フレーム

（一定時間）毎のパワーを計算するパワー計算部11aと、1フレーム毎の1次および7次の自己相関係数を算出する自己相関係数算出部11bと、1フレーム毎の1次および3次のケプストラム係数を算出するケプストラム係数算出部11cとから構成される。これらの特徴量は入力信号の母音性を検出するために用いられる。

【0012】次に、12はあらかじめ多数の母音の学習データについて特徴抽出部11で抽出した特徴量を用いて母音毎の平均値と共分散行列を算出し、母音毎の標準モデルを作成する母音標準モデル作成部である。13は特徴抽出部11から出力されるフレーム毎の入力信号の1次および7次の自己相関係数と1次および3次のケプストラム係数について、母音標準モデル作成部12にて作成した各母音標準モデルとの対数尤度を計算する尤度計算部であり、14は尤度計算部13にて計算に用いたフレームの前後数フレームにおいて、尤度計算部13で同様に計算された対数尤度を用いて、各母音毎にフレーム平均対数尤度を計算し、ある適当な閾値と比較することでその入力信号数フレームが母音であるかどうかを判定する母音判定部である。15はパワーの一定レベル以上の入力信号の塊について母音判定部14によりいずれかの母音と判定されたフレームの個数がある適当な閾値以上のときにその塊を音声と判定する最終判定部である。

【0013】以下、本発明の一実施例について（図1）のブロック構成図を参照しながら詳細に説明する。音響信号がマイクロホンを通して入力されると、特徴抽出部11でまず複数の特徴量が抽出される。パワー計算部11aでは、一定時間毎のパワー値が例えば（数1）で算出される。一定の時間間隔は、ここでは例えばサンプリング周波数を10KHzとして、200点（20ms）とし、この時間単位をフレームと呼ぶ。

【0014】

【数1】

$$P_i = \sum_{k=1}^{200} S_k * S_k$$

【0015】ここで、 P_i はフレーム*i*でのパワー値、 S_k はフレーム内の入力信号のサンプル値を示す。このパワー値は発声条件の違いによるパワーの違いを統一して扱えるように、パワーの大きな区間内の最大値、最小値間を例えば0から1までの値に正規化して用いる。自己相関係数算出部11bではフレーム毎に1次の自己相関係数 $A_i(1)$ が（数2）、7次の自己相関係数 $A_i(7)$ が、（数3）で算出される。

【0016】

【数2】

$$A_i(1) = \sum_{k=1}^{200} S_k * S_{k+1}$$

【0017】

【数3】

$$A_i(7) = \sum_{k=1}^5 S_k * S_{k+1}$$

【0018】さらに $A_i(1)$ 、 $A_i(7)$ は0次の自己相関係数で正規化される。ケプストラム係数算出部11cでは、フレームiでの1次および3次のケプストラム係数 $C_i(1)$ 、 $C_i(3)$ が線形予測分析により求められる。

【0019】母音標準モデル作成部12では、あらかじめ多数の音声データの母音部分について特徴抽出部11で得られる特徴量を抽出しておき、これらの特徴量を用いて次の方法により母音毎の平均値と共分散行列を算出し、母音毎の標準モデルを作成する。すなわち、母音データとしては母音kの学習用データ y_N （データ数N）を用い、 y_N がm次元の多次元正規分布に従うと仮定できる場合、その平均値 μ_k と共分散行列 Σ_k を（数4）、（数5）のように計算にて求めることができる。

【0020】

【数4】

$$\mu_k = \frac{1}{N} \sum_{n=1}^N y_n$$

【0021】

【数5】

$$\Sigma_k = \frac{1}{N} \sum_{n=1}^N (y_n - \mu_k)(y_n - \mu_k)'$$

*

$$L_{ik} = -\frac{1}{2} (x_i - \mu_k)' \Sigma_k^{-1} (x_i - \mu_k) - \frac{1}{2} \ln |\Sigma_k| + C$$

【0025】ただし、 x_i はm次元のベクトル（m次元の特徴パラメータ）であり、tは転値、 -1 は逆行列、Cは定数を示す。

【0026】母音判定部14は、母音の時間的な継続性を表現するため、検出しようとする目的フレームの前後Nフレーム（これをセグメントという）を用いて、母音判定を行う。各母音毎に尤度計算部13にて計算された対数尤度 L_{ik} を用いて、次の条件式（数7）を満たせばそのセグメントは母音であるとみなす。

【0027】

【数7】

$$\frac{1}{2 * N + 1} \sum_{t=-N}^N L_{ik} \geq L_{kTH}$$

【0028】ただし、 L_{kTH} は母音標準モデルkに関する判別閾値（フレーム平均対数尤度の閾値）である。

【0029】このように、各特徴パラメータの影響を効果的に、しかも総合的に判定できる評価値（対数尤度）を用いることで、各特徴パラメータ毎に閾値を設定する方法よりも、定常雑音が付加された場合などのような特徴量の値の変動に対して頑健なシステムが構築できる。また、多くの閾値をヒューリスティックな方法により決

6

*【0022】ただし、tは転値を示す。これにより、母音毎の標準モデルのモデル形状（平均値 μ_k 、及び分散 Σ_k ）が求められる。ただし、 y_N 、 μ_k はm次元のベクトル（m次元の特徴パラメータ）であり、 Σ_k は $m \times m$ 次元のマトリックスである。母音データとしては例えば、ある標準話者の母音kの学習用データとして母音部分を切り出し、母音中心フレーム±2フレームのデータを用いればよい。また、複数の話者のデータを用いることで、話者の発声の変動に強い標準パターンを作成することができる。

【0023】尤度計算部13は、特徴抽出部11から出力されるフレーム毎の入力信号のいくつかの特徴パラメータについて、母音標準モデル作成部12にて作成した各母音標準モデルとの対数尤度を計算する部分である。母音検出に用いる距離尺度は、使用する各特徴パラメータの分布を多次元正規分布と仮定した場合の統計的距離尺度である。母音毎の標準モデルkに対する、iフレーム目の入力ベクトル（スペクトル） x_i の対数尤度 L_{ik} は、（数6）で計算される。

20 【0024】

【数6】

定する必要がある利点がある。さらに、音響信号数フレーム分を1塊に考えて判定することで、母音などのような継続的な音声に対してより有効な判定法となっている。

【0030】最終判定部15は、パワーの一定レベル以上の入力信号の塊についての母音サンプルの数がある適当なしきい値以上のときにその塊を音声と判定する最終判定部である。最終判定部15では、まずパワー計算部11aで得られたパワー値系列からあらかじめ定めたパワーしきい値を決められた長さ以上越える区間を音声候補区間として検出する。この音声候補区間内において、母音判定部14により母音kと判定されたセグメントの個数を数え、母音kと判定された母音セグメントの数を C_k 、あらかじめ定めた区間内の母音セグメント数のしきい値 M_k とすると、（数8）の条件を満たすならば、この音声候補区間は音声であると判定する。これを全ての母音について行う。

【0031】

【数8】

$$C_k \geq M_k$$

【0032】以下に、実際に本方法により実験した結果を示す。（表1）に、本手法で用いた4つの特徴パラメ

ータを示す。これらの特徴パラメータは予備実験の結果、音声と他の非定常雑音との分離が比較的良く、またLPCケプストラム係数の算出過程において容易に得られるパラメータである。まず、1次の正規化自己相関係数及び1次の線形予測係数は有声／無声の判別に適したパラメータであり、7次の正規化自己相関係数は低周波性の雑音を区別するのに適したパラメータである。また、3次のLPCケプストラム係数は、5母音の中でも／i／に特徴的な性質を示すパラメータである。

【0033】

【表1】

- ・ 1 次 の 自 己 相 関 係 数
- ・ 1 次 の 線 形 予 測 係 数
- ・ 7 次 の 自 己 相 関 係 数
- ・ 3 次 の L P C ケ プ ス ト ラ ム 係 数

*

非 定 常 雑 音 グ ル ー プ	サ ン プ ル 数
紙 の 音	2 8 5
手 や 物 で 机 を 叩 く 音	2 1 5
ガ ラ ス ・ 食 器 の 触 れ 合 う 音	1 0 0
マイクに触れる・叩く音	1 8 4
息 吹 き ・ 咳 払い	1 1 3

【0036】

【表3】

サンプリング周波数	10[kHz]
分析窓	20[ms]ハミング窓
フレーム周期	10[ms]
プリアンファシス	1-0.9z ⁻¹
LPC分析次数	12次

【0037】男性話者5名分の母音データから標準モデルの作成を行い、標準話者を含む10名の話者についての音声検出実験及び（表2）の非定常雑音の除去実験を行った。（図2）は、母音セグメント長を1フレームから11フレームまで変化させたときの、音声検出率と雑音誤検出率の関係を示したものである。判別閾値を適当に変化させることで、検出性能の最適値を求めることができるが、5フレーム以上ではほとんど判別性能に差はない。結局、母音セグメント長7フレームで判別閾値=-1.2のとき、音声検出率99.3%（雑音誤検出率9.0%）が得られた。

【0038】次に、本手法の定常騒音下での検出性能を評価するために、白色雑音を付加したときのS/N比と検出率との関係を調べた。（図3）は、母音セグメント長を7フレームに固定したときの、各S/N比に対する音声検出率と雑音誤検出率との関係を示したものである。その結果、検出性能はS/N比が12dBまでほと

*【0034】音声データは、男性10名の発声した日本語200単語である。標準モデルの作成には、標準話者の発声した各母音の音韻中心±2フレームを使用した。但し、計算効率を考慮して、各パラメータ間の相関はないとし、共分散行列の対角成分のみを計算に用いた。雑音データとしては、（表2）に示す5雑音グループ（約900サンプル）の非定常雑音を用いた。また、分析条件を（表3）に示す。

【0035】

10 【表2】

んど影響を受けていない。

【0039】以上のように本実施例の音声検出装置によれば、入力信号から一定時間毎の音声の複数の特徴量を抽出する特徴量抽出部11と、あらかじめ多数の母音に関する学習データについてフレーム単位で抽出した前記特徴量を用いて母音毎の平均値と共分散行列を算出し、母音毎の標準モデルを作成する母音標準モデル作成部12と、入力信号から得られた複数の音声の特徴量と母音標準モデル作成部にて作成した各母音標準モデルとの対数尤度を計算する尤度計算部13と、母音判定を行うフレームの前後数フレーム分の対数尤度を用いて、各母音毎にフレーム平均対数尤度を計算し、ある適当なしきい値とを比較することでその入力信号数フレーム分が母音であるかどうかを判定する母音判定部14と、パワーの一定レベル以上の入力信号の塊について母音判定部14によりいずれかの母音と判定された母音サンプルの個数がある適当なしきい値以上のときにその塊を音声と判定

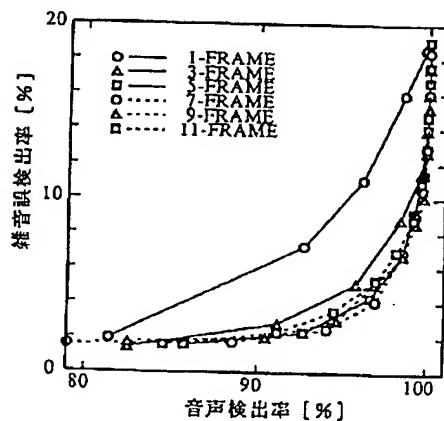
する最終判定部15とを具備して構成することにより、比較的簡単な構成で様々な音響信号の中の音声を正確に判定することができる音声検出装置を提供することができる。

【0040】なお、上記の実施例においては、特徴抽出部において入力信号の母音性を検出するための特徴量として自己相関係数とケプストラム係数を用いた例で説明したが、これに限定されず、偏自己相関関数やメルケプストラム係数などを用いてもかまわない。

【0041】

【発明の効果】以上の実施例から明らかなように本発明によれば、音声の特徴付ける複数の特徴量を抽出し、多数の学習用母音データを用いて母音標準モデルを作成しておき、入力信号から得られた複数の特徴量と母音標準モデルから得られる対数尤度を計算し、数フレーム分を一括して母音の検出を行って音声を検出するように構成しているので、比較的簡単な構成で入力信号が音声かそれ以外かを正確に判定することができる音声検出装置を

【図2】



提供することができる。

【図面の簡単な説明】

【図1】本発明の一実施例の音声検出装置の全体構成を示すブロック図

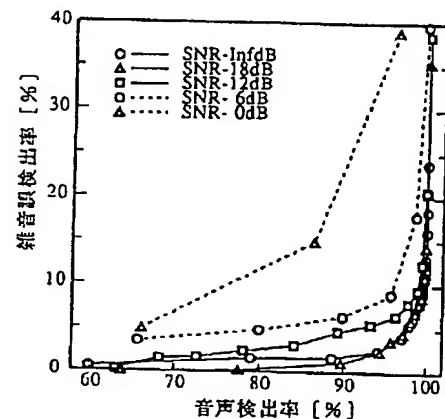
【図2】母音セグメント長の影響による音声検出率と雑音誤検出率の関係を示す図

【図3】S/N比の影響による音声検出率と雑音誤検出率の関係を示す図

【符号の説明】

- 10 11 特徴抽出部
11 a パワー算出部
11 b 自己相関係数算出部
11 c ケプストラム係数算出部
12 母音標準モデル作成部
13 尤度計算部
14 母音判定部
15 最終判定部

【図3】



【図1】

11. 特徴抽出部

